

Spatial models for point and areal data using Markov random fields on a fine grid

Christopher J. Paciorek

Department of Biostatistics, Harvard School of Public Health

and

Department of Statistics, University of California, Berkeley

Abstract

I consider the use of Markov random fields (MRFs) on a fine grid to represent latent spatial processes when modeling point-level and areal data, including situations with spatial misalignment. Point observations are related to the grid cell in which they reside, while areal observations are related to the (approximate) integral over the latent process within the area of interest. I consider several approaches to specifying the neighborhood structure for constructing the MRF precision matrix, presenting results comparing these MRF representations analytically, in simulations, and in two examples. I highlight the use of an MRF that approximates a thin plate spline as an alternative to standard CAR models. In particular, I show that, in contrast to this thin plate spline approximation, various neighborhood weighting approaches presented in the literature do not produce smooth fields, even with higher-order neighbors.

1 Introduction

Markov random field (MRF) models (also called conditional autoregressive (CAR) models) are the dominant approach to analyzing areally-aggregated spatial data, such as disease counts in adminis-

trative units. In contrast, point-referenced data are generally modeled using Gaussian process (GP) models that posit a continuous latent underlying spatial surface. The popular approach of kriging for spatial prediction can be seen as a GP model when normality is assumed.

GP models are computationally challenging for large datasets because of manipulations involving large covariance matrices, and there has been a large body of recent work attempting to reduce the computational burden through reduced rank approximations (Kammann and Wand, 2003; Banerjee, Gelfand, Finley, and Sang, Banerjee et al.; Sang and Huang, 2012), inducing sparseness in the covariance (covariance tapering) (Furrer et al., 2006; Kaufman et al., 2008; Sang and Huang, 2012), approximate likelihood approaches (Stein et al., 2004), and approximation in the Fourier domain (Wikle, 2002; Paciorek, 2007), among others. In contrast, MRF models work with the precision matrix directly, so calculation of the likelihood is computationally simple. In MCMC implementations, one can exploit the Markovian structure to sample the value of the field for each area sequentially or exploit the sparsity of the precision matrix when sampling the field values jointly (Rue and Held, 2005).

Given the computational attractiveness of the MRF approach, it is appealing to consider its use for point-referenced data as well. One issue lies in how to define the neighborhood structure for a set of point-referenced observations, but a second more fundamental issue lies in the fact that the model structure changes when one changes the number of sites under consideration. The solution that I highlight here is to relate the point-referenced observations to an underlying regularly-gridded surface modeled as an MRF, which approximates a smooth surface using a fine piecewise approximation.

Considering an underlying gridded surface for areally-aggregated data is also appealing. Kelsall and Wakefield (2002) argue for the use of a smooth underlying surface to model areal data, with each areal observation related to the average of the surface over the defined area. They found approximations to the integrals involved and worked with an underlying GP representation; related work includes Fuentes and Raftery (2005) and Hund et al. (2012). Here I relate areal observations to an underlying MRF on a fine grid, approximating the necessary integrals as a simple weighted

average of the MRF values for the grid cells that overlap each area.

The approach of building a model for areal data based on a smooth underlying surface is also appealing as a way to deal with spatial misalignment in different data sources, the so-called modifiable areal unit problem in geography. Mugglin et al. (2000) suggest that one work with intersections of different grids and build a coherent model at the resolution of the intersections. Instead, I propose that one relate all observations, whether point observations or areal observations that may be spatially misaligned, to a common underlying grid to produce an aggregation-consistent model, as done in Paciorek (2012) using the methodology discussed here. The strategy of relating either point or irregular area observations to an MRF on a fine grid was also suggested in Besag and Mondal (2005).

The most common form of MRF represents the spatial dependence structure such that areas that share a boundary are considered neighbors, with an area conditionally independent (given its neighbors) of any non-neighboring areas, a so-called first-order neighborhood structure. Fig. 1 shows the results of using this MRF structure on a fine regular grid to model point-referenced data. The fitted smooth surface is not visually appealing, in contrast to an MRF that approximates a thin plate spline (Rue and Held, 2005) and to kriging. The results are consistent with Besag and Mondal (2005), who show that the intrinsic (i.e., improper) first-order MRF on a two-dimensional regular grid produces spatial fields whose distribution approaches two-dimensional Brownian motion (the de Wijs process) asymptotically as the grid resolution increases. Given this continuous but non-differentiable representation of the underlying surface, the local heterogeneity of the surface estimate in Fig. 1 is not surprising. However, note that Besag and Mondal (2005) and Besag in his comments on Diggle et al. (2010) argue that the de Wijs process is preferable to GPs in the Matérn class.

A common alternative to this standard nearest-neighbor structure is to extend the neighborhood structure beyond bordering areas (Pettitt et al., 2002; Hrafnkelsson and Cressie, 2003; Song et al., 2008). Such higher-order neighborhood structure might be expected to produce more smooth process representations, but I show that straightforward higher-order neighborhoods do not achieve

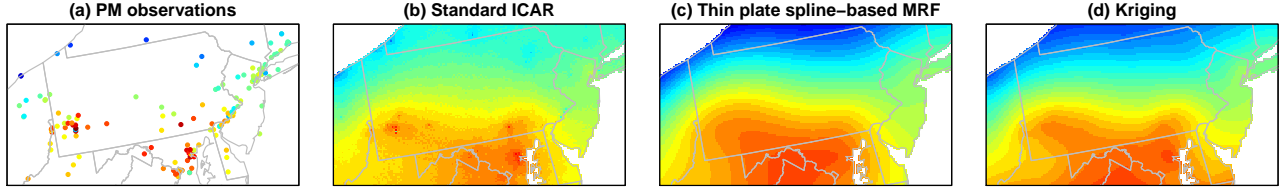


Figure 1: Example fits to particulate matter (PM) air pollution point observations (a) using the standard first-order MRF on a fine regular grid (b), an MRF approximation to a thin plate spline, again represented on a fine regular grid (c), and using kriging with an exponential covariance (d).

this. An alternative that I highlight here is an MRF approximation to a thin plate spline (TPS) that involves only nearby grid cells as neighbors (Rue and Held, 2005; Yue and Speckman, 2010). Finally, Lindgren et al. (2011) have recently developed a powerful theory and methodology for approximating GPs in the Matérn class that will likely be widely used. The thin plate spline approximation is a limiting case of the Lindgren et al. (2011) representation, and my results shed light on the distinctions between standard first-order MRF models, the thin plate spline approximation, and GP representations.

In this paper, I present a general model for both areal and point-referenced data that deals simply with spatial misalignment by using an MRF on a fine regular grid. In light of the lack of smoothness of the popular first-order MRF seen in Fig. 1, I compare several approaches for the MRF neighborhood structure using both analytic calculations and simulations. The results suggest that the thin plate spline approximation often performs well and help to explain the strange behavior of the first-order MRF. The results also suggest that higher-order neighborhood structures can have unappealing properties when they are not constructed as approximations to particular spatial process representations. I discuss computation, focusing on the difficulties in fitting non-normal data, and considering the use of PQL as well as the widely-used INLA approximation (Rue et al., 2009). I present two examples, one showing the use of the approach for point data and the second with areal data.

2 Spatial model for point and areal data

2.1 Model structure

Here I present a basic model for exponential family data. Let $\mu_i = E(Y_i | \mathbf{X}_i, \mathbf{g})$ be related via a link function, $h(\cdot)$, to a linear predictor:

$$h(\mu_i) = \mathbf{X}_i^\top \boldsymbol{\beta} + \mathbf{K}_i \mathbf{g}, \quad (1)$$

where \mathbf{K}_i is the i th row of a mapping matrix, \mathbf{K} , discussed further below. I represent the unknown, latent spatial process, $g(\cdot)$, as a piecewise constant surface on a fine rectangular grid,

$$\mathbf{g} \sim \mathcal{N}(\mathbf{0}, (\kappa \mathbf{Q})^-) \quad (2)$$

where \mathbf{Q} is an MRF precision matrix and κ a precision parameter, recognizing the potential singularity of \mathbf{Q} by using the generalized inverse notation. Computational issues related to the singularity are discussed in Section 4.

For a point-referenced datum, \mathbf{K}_i will be a sparse vector with a single 1 that matches the location of the observation to the grid cell in which it falls. Note that one may include covariates in \mathbf{X} that can help to account for within grid cell heterogeneity. For an areally-aggregated datum, I consider the relevant functional of the surface to be the average value of the underlying surface over the area A_i : $g^*(A_i) \equiv \int_{A_i} g(\mathbf{s}) d\mathbf{s}$. Using the piecewise representation, this can be approximated as $g^*(A_i) \approx \sum_{j \in A_i} w_j g_j$ where $j \in A_i$ indexes the grid cells that the A_i overlaps and w_j is the proportion of A_i that falls in the j th grid cell. Hence the non-zero elements of \mathbf{K}_i contain the proportions, $\{w_j\}$, as the appropriate elements. If desired, one could interpolate between the values of \mathbf{g} at the grid cell centroids for a more accurate approximation.

2.2 Potential MRF models

Here I present the MRF models that I consider for \mathbf{g} and their corresponding precision matrices, \mathbf{Q} . I consider only intrinsic models with singular precision matrices. These specify improper priors with respect to one or more linear combinations of the process values, as the eigenvectors of \mathbf{Q} with zero eigenvalues have infinite prior variance (Banerjee et al., 2003). My focus on intrinsic models is motivated by noting that proper first-order MRF models tend not to allow high correlations between neighbors unless the precision matrices are close to singular (Banerjee et al., 2003; Wall, 2004). Furthermore, intrinsic models represent the conditional mean for the process in a given area as a weighted average of the process values of the neighboring areas with weights summing one, while proper models have weights summing to less than one.

1. Traditional intrinsic CAR model (ICAR): The most commonly-used MRF model is a simple first-order model that treats any two areas that share a border as neighbors. The corresponding precision matrix has diagonal elements, Q_{ii} , equal to the number of neighbors for the i th area, while $Q_{ij} = -1$ (the negative of a weight of one) when areas i and j are neighbors and $Q_{ij} = 0$ when they are not. On a grid, the simplest version of this model treats the four nearest grid cells as neighbors (i.e., cardinal neighbors). Besag and Mondal (2005) show that the model is asymptotically equivalent to two-dimensional Brownian motion, and Lindgren et al. (2011) show that this model approximates a GP in the Matérn class (3) with the spatial range parameter $\rho \rightarrow \infty$ and differentiability parameter $\nu \rightarrow 0$.
2. Extended neighborhood model (DICAR and HICAR): One might generalize the first-order Markovian structure of the ICAR model to allow for higher-order dependence by considering areas that do not share a border (but are close in some sense) to be neighbors. At its simplest, this simply introduces additional values of -1 off the diagonal of \mathbf{Q} as in Song et al. (2008). I will call this model the higher-order ICAR (HICAR). A more nuanced version would have the weight for a pair of areas depend on the distance between the two areas (usually declining with distance), such as Euclidean distance or the number of intervening cells between the

two areas (Pettitt et al., 2002; Hrafnkelsson and Cressie, 2003). Then $Q_{ij} = -\delta(i, j)$ where $\delta(\cdot, \cdot)$ is the chosen weight or distance function. I will term this model the distance-based ICAR (DICAR) and implement the model using the function in Hrafnkelsson and Cressie (2003), $\delta(i, j) = -d_{ij}^{\log .05 / \log r}$, where d_{ij} is the distance between area centroids and r is the (user-chosen) maximum distance at which the weight is non-zero.

3. Thin plate spline approximation (TPS-MRF): Rue and Held (2005, Sec. 3.4.2) present a second-order model where the weights on different order neighbors are constructed so the model approximates a thin plate spline. The nearest cardinal neighbors have a weight of 8 ($Q_{ij} = -8$), the nearest diagonal neighbors a weight of -2 ($Q_{ij} = 2$) and the second nearest cardinal neighbors a weight of -1 ($Q_{ij} = 1$). Note the presence of negative weights, unlike in most MRF models with higher order neighborhood structures. Paciorek and Liu (2012, App. C) describe the derivation of the full \mathbf{Q} matrix, including boundary effects, in detail. In one dimension, this model corresponds to the widely-used second-order auto-regressive model (an IID model on second differences) (Breslow and Clayton, 1993). Lindgren et al. (2011) show that this model approximates a GP in the Matérn class (3) with $\nu = 1$ and the spatial range parameter $\rho \rightarrow \infty$.

In this work I compare MRF models to GPs in the Matérn class, where I parameterize the Matérn correlation function as

$$R(d) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{2\sqrt{\nu}d}{\rho} \right)^\nu \mathcal{K}_\nu \left(\frac{2\sqrt{\nu}d}{\rho} \right), \quad (3)$$

where d is Euclidean distance, ρ is the spatial range parameter, and $\mathcal{K}_\nu(\cdot)$ is the modified Bessel function of the second kind, whose order is the smoothness (differentiability) parameter, $\nu > 0$. $\nu = 0.5$ gives the exponential covariance.

3 Comparing MRF Structures

3.1 Covariance and smoothing properties

3.1.1 Eigenstructure

To better understand the spatial dependence structures implied by the various MRF precision matrices, I quantify the magnitude of the variability for different scales (frequencies) of spatial variability by considering the eigendecomposition of \mathbf{Q} . Given the model $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^-)$, we can consider the eigendecomposition, $\mathbf{Q}^- = \mathbf{\Gamma} \mathbf{\Lambda}^- \mathbf{\Gamma}^\top$. To generate realizations of \mathbf{g} , we have $\mathbf{g} = \mathbf{\Gamma} \mathbf{u}$ for $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}^-)$. Thus the inverse eigenvalues quantify the magnitude of variability associated with patterns or modes of variability represented by the eigenvectors.

Empirical exploration indicates that the inverse eigenvalues decline as the frequency of variability represented in the eigenvectors increases. Therefore to visualize how different MRF models weight variability at different frequencies, I plot the ordered inverse eigenvalues. Because the matrix \mathbf{Q} is multiplied by a scalar precision, only the relative magnitudes are of interest, so I scale the inverse eigenvalues such that the 100th largest inverse eigenvalue for each precision matrix is taken to be equal to one. While the eigenvectors for the various precision matrices are not the same, empirical exploration shows that they are quite similar and represent very similar spatial scales of variation for a given position in the ordering of the eigenvector/value pairs. Furthermore, projection onto a common set of eigenvectors (from the ICAR model) confirms that the slight differences in eigenvectors between models do not impact the results shown next.

In the comparisons, I compare the TPS-MRF model to a GP with Matérn correlation with $\nu = 2$, rather than $\nu = 1$ as would be natural given the Lindgren et al. (2011) relationship. I do this because thin plate splines represent smooth functions and the Matérn has M mean square derivatives when $\nu > M$ (Stein, 1999, p. 32). I compare the ICAR model to a GP with an exponential correlation (Matérn with $\nu = 0.5$) despite the Lindgren et al. (2011) result relating the ICAR model to a Matérn covariance with $\nu \rightarrow 0$ because the Matérn covariance is only valid for $\nu > 0$.

Fig. 2 plots the size of the ordered inverse eigenvalues for different MRF precision matrices and in comparison to GP models for a regular spatial grid of dimension 75×75 . The TPS-MRF puts more weight on the lower frequency eigenvectors and less weight on the higher frequency eigenvectors than the ICAR. This is not surprising given the relationship of the ICAR model to Brownian motion and the smoothness of splines. The TPS-MRF eigenvalue curve lies within the set of eigenvalue curves (for varying values of the range parameter, ρ) from the Matérn model with $\nu = 2$. At low frequency, the ICAR eigenvalue curve lies within the set of eigenvalue curves from the exponential model, while the ICAR model puts more weight on high frequency eigenvectors than the exponential model, consistent with the ICAR approximating a Matérn covariance with $\nu \rightarrow 0$ (Lindgren et al., 2011).

Particularly interesting is the behavior of the higher-order and distance-based ICAR variations, which behave similarly to the ICAR model (Fig. 2b). At low frequency, the curves coincide, while the HICAR and DICAR models put more weight on the higher frequencies than the ICAR model. This indicates that one cannot use these approaches to represent surfaces smoother than those represented by the ICAR model. In fact, based on this analysis, it is unclear why one would use these representations, given that the motivation for their use has been to provide more smoothness than the ICAR model.

3.1.2 Equivalent kernels

To understand the smoothing behavior of the various models, I next consider their equivalent kernels, which quantify the local averaging that the models do to make predictions. Under a normal likelihood and ignoring covariates for simplicity, $\mathbf{Y} \sim \mathcal{N}(\mathbf{K}\mathbf{g}, \tau^2\mathbf{I})$. When there are observations at all the grid cells, the smoothing matrix, \mathbf{S} , in $\hat{\mathbf{g}} = \mathbf{S}\mathbf{y}$ can be expressed as

$$\mathbf{S} = \frac{1}{\tau^2} \left(\kappa \mathbf{Q} + \frac{1}{\tau^2} \mathbf{I} \right)^{-1} = (\lambda \mathbf{Q} + \mathbf{I})^{-1},$$

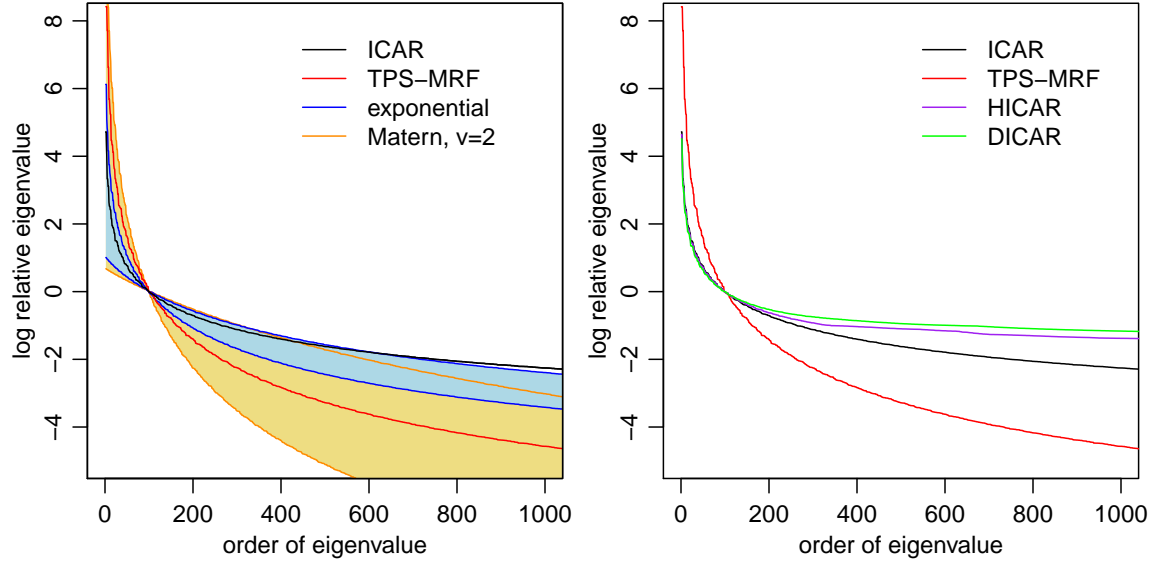


Figure 2: (a) Log inverse eigenvalues (relative to the 100th largest inverse eigenvalue) for the ICAR and TPS-MRF models, in comparison with the log inverse eigenvalues for the Matérn covariance (with $\nu = 2$) and exponential covariance models. The yellow-orange shading, bounded by orange lines, indicates the range of eigenvalue curves obtained for the Matérn as the range parameter, ρ , varies from the size of the full domain in one dimension to $1/25$ of the domain in one dimension. Similarly, the blue shading, bounded by dark blue lines, is for the exponential covariance as the range varies from the full domain to $1/25$ of the domain. (b) Log inverse eigenvalues (relative to the 100th largest inverse eigenvalue) for the ICAR and TPS-MRF models, both repeated from (a); the HICAR model with grid cells up to and including three units in distance considered neighbors; and the DICAR model with non-zero weights that decay with distance up to and including five units in distance. Only the first 1000 inverse eigenvalues are plotted, as by the 1000th, the features represented in the eigenvectors occur within groups of 5-10 grid cells, near the limit of scales that could be resolved in a gridded representation.

where $\lambda \equiv \tau^2\kappa$. For the GP models, we have $\hat{\mathbf{g}} = \hat{\mu}\mathbf{1} + \sigma^2\mathbf{R}_\theta(\sigma^2\mathbf{R}_\theta + \tau^2\mathbf{I})^{-1}(\mathbf{y} - \hat{\mu}\mathbf{1})$ so

$$\mathbf{S} = \sigma^2\mathbf{R}_\theta(\sigma^2\mathbf{R}_\theta + \tau^2\mathbf{I})^{-1} = (\lambda\mathbf{R}_\theta^{-1} + \mathbf{I})^{-1},$$

where τ^2 is the error variance, σ^2 is the marginal variance of the GP, and \mathbf{R}_θ is the correlation matrix, a function of parameter(s), θ . $\lambda \equiv \tau^2/\sigma^2$ can be thought of as a smoothing parameter, as in the MRF model. Note that the sum of the weights is not one when expressing \mathbf{S} as above, because it ignores $\hat{\mu}$, which is also linear in the observations. I ignore this component of the prediction as it involves adding and subtracting a constant that does not vary by location.

Fig. 3 shows the various kernels, plotted as a function of one of the cardinal directions, with the other direction held fixed, for two values of the smoothing parameter, λ . Fig. 3a,d shows that the ICAR kernel puts less weight near the focal cell and more further from the focal cell, compared to the TPS-MRF, but with a spike at the focal cell. This behavior helps to explain the bulls-eye effect and the greater shrinkage towards an overall mean in the gaps between observations seen for the ICAR model in Fig. 1. The TPS-MRF kernel shows some small-magnitude negative weights, which is consistent with the negative weights in the equivalent kernels for spline smoothing and Gaussian process smoothing (Silverman, 1984; Sollich and Williams, 2005). In Fig. 3b,e, we see that the HICAR (with neighbors within three units) and DICAR (with non-zero weights within five units) models place very little weight near the focal observation and spread their weight very widely. The result is that their kernels are even more extreme than the ICAR kernel in making predictions that heavily weight observations far from the focal location. Results are similar but not as extreme for HICAR and DICAR models with smaller neighborhoods. As with the eigenvector analysis, this suggests the HICAR and DICAR models have little practical appeal. Fig. 3c,f compares the ICAR and TPS-MRF to the equivalent kernels for GP models with ρ set to one-tenth of the domain in one dimension. The ICAR model shows some similarity to the exponential-based GP in terms of tail behavior and the spike at the focal cell. The TPS-MRF equivalent kernels are rather different than the Matérn -based GP model with $\nu = 2$, but more similar in terms of tail

behavior than when compared to the exponential-based model.

Fig. 4 reinforces these points, showing image plots of the equivalent kernels in two dimensions for the ICAR and TPS-MRF models, as well as the Matérn and exponential covariance models. The ICAR kernel puts little weight near the focal cell but spreads positive weight further from the focal cell than the other approaches. The exponential is qualitatively similar, but less extreme than the ICAR. The TPS-MRF and Matérn models are somewhat similar, but the TPS-MRF puts more weight near the focal cells.

3.2 Analytic assessment of predictive ability

For this analysis, I derive the expected mean squared error, averaging over randomness in observations and randomness in latent spatial process realizations. I assume the simple generative model

$$\begin{aligned} \mathbf{Y} &\sim \mathcal{N}(\mathbf{K}\mathbf{g}, \tau^2 \mathbf{I}) \\ \mathbf{g} &\sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{C}), \end{aligned} \tag{4}$$

with $\mathbf{C} = \sigma^2 \mathbf{R}_\theta$, where \mathbf{R}_θ is a correlation matrix based on either the exponential or Matérn with $\nu = 2$. Here \mathbf{X} represents spatial basis functions, including an overall mean and possibly linear and higher-order functions of the spatial coordinates, as in ordinary kriging and universal kriging. I consider six scenarios with the domain the unit square: $n = 100$ with locations uniformly sampled, $n = 100$ with observations sampled according to a Poisson cluster process (PCP) with parent intensity of 25 and standard normal kernels with standard deviation of 0.05, $n = 1000$ with uniform and PCP sampling, $n = 10000$ on a regular 100×100 grid, and $n = 100$ areal observations on a coarse 10×10 grid overlaid on a regular 100×100 grid. For each scenario, I use a full factorial design with respect to $\nu \in \{0.5, 2\}$, $\rho \in \{0.005, 0.02, 0.08, 0.32, 1.28, 2.56\}$, and $\tau^2 \in \{0.05^2, 0.15^2, 0.45^2, 1.35^2\}$. For $n = 100$ and $n = 1000$, I use 10 replicates to average over the random sampling of locations; the resulting Monte Carlo standard errors are small relative

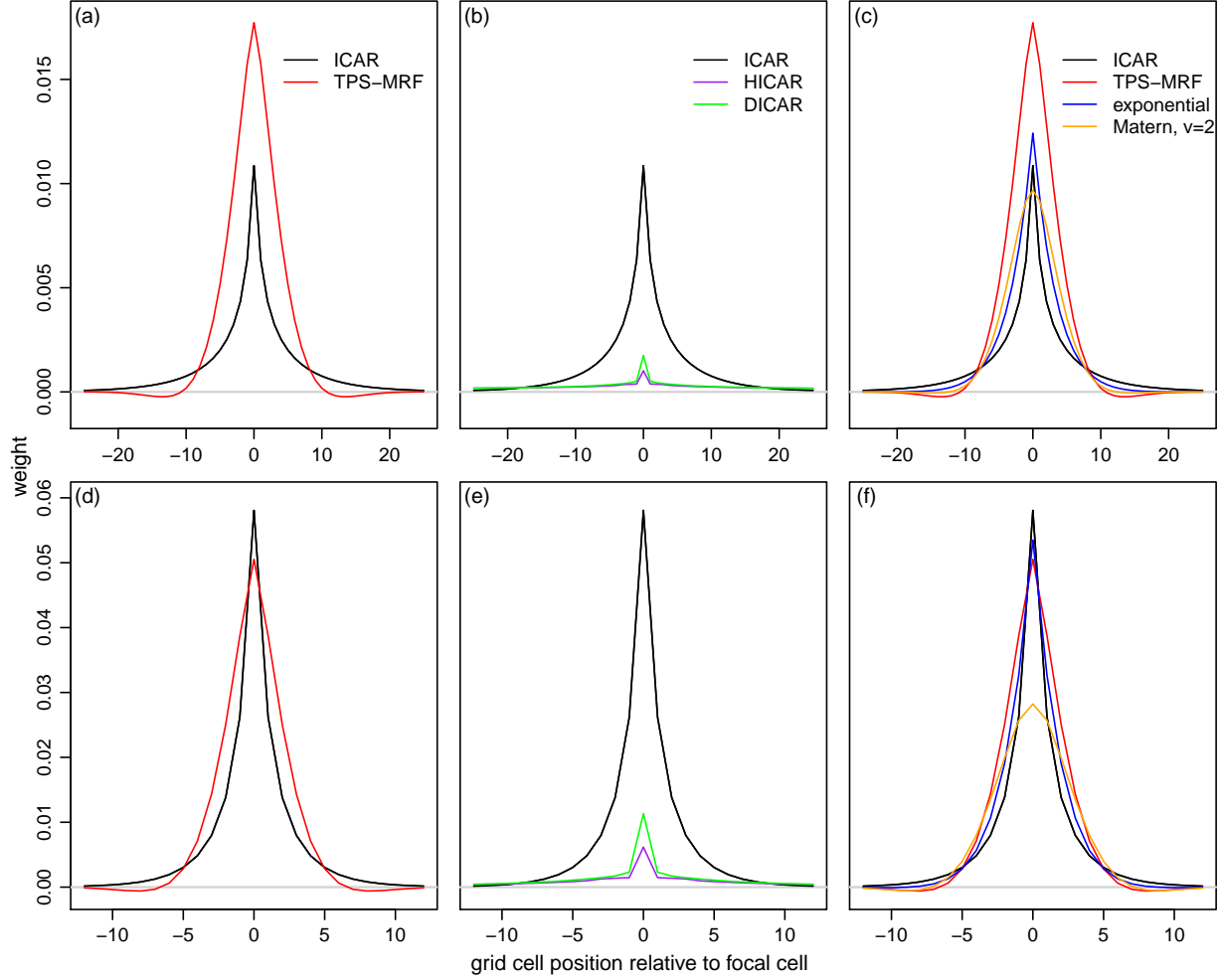


Figure 3: Equivalent kernel cross-sections in one of the cardinal directions for various models with (a-c) $\lambda = \exp(4)$, which produces less localized weighting and (d-f) $\lambda = \exp(2)$, which produces more localized weighting. In (b) and (e), the HICAR model has grid cells up to and including three units in distance considered neighbors, and the DICAR model has decaying non-zero weights up to and including five units in distance. In (c) and (f), the Matérn (with $\nu = 2$) and exponential covariance models have range, ρ , set to one-tenth of the domain size in one dimension.

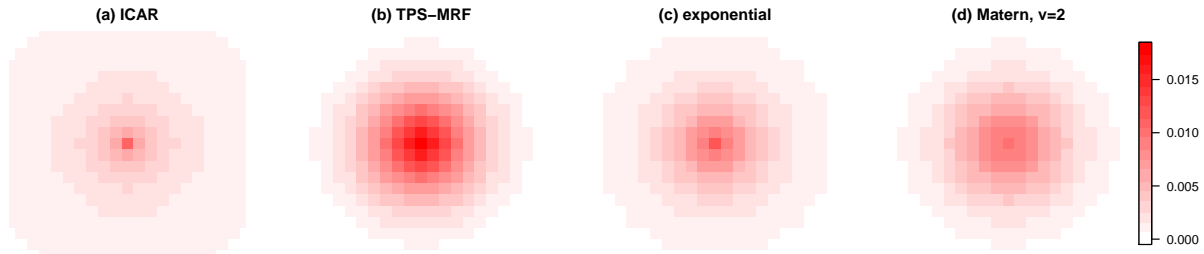


Figure 4: Equivalent kernels in two dimensions for $\lambda = \exp(4)$: (a) ICAR, (b) TPS-MRF, (c) exponential covariance with range, ρ , equal to one-tenth of the domain in one dimension, and (d) Matérn covariance with $\nu = 2$ and that same range.

to the point estimates of the reported MSE values. Data and predictions are taken to occur at the grid cell centroids, but multiple observations are allowed to occur in a grid cell. For simplicity, I simulate with $\sigma^2 = 1$ and $\beta = \mathbf{0}$.

I consider fitting the data using either the ICAR or TPS-MRF models on a regular 100×100 grid:

$$\mathbf{g} \sim \mathcal{N}(\mathbf{X}\beta, (\kappa\mathbf{Q})^-),$$

where $\mathbf{X}\beta$ represents the overall mean in the ICAR model and the mean plus linear functions of the spatial coordinates in the TPS-MRF model. One can express the best prediction for \mathbf{g} as

$$\begin{aligned} \hat{\mathbf{g}} &= E(\mathbf{g}|\mathbf{Y}, \cdot) = (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{Q})^{-1}(\mathbf{K}^\top \mathbf{Y} + \lambda \mathbf{Q} \mathbf{X} \hat{\beta}) \\ &= (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{Q})^{-1}(\mathbf{K}^\top + \lambda \mathbf{Q} \mathbf{X} ((\mathbf{K} \mathbf{X})^\top \Sigma^{*-1} \mathbf{K} \mathbf{X})^{-1} (\mathbf{K} \mathbf{X})^\top \Sigma^{*-1}) \mathbf{Y} \\ &= (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{Q})^{-1} \mathbf{K}^\top \mathbf{Y} \\ &\equiv \mathbf{S} \mathbf{Y}. \end{aligned} \tag{5}$$

where $\Sigma^* \equiv \mathbf{I} + \mathbf{K}(\lambda \mathbf{Q})^- \mathbf{K}^\top$, $\lambda \equiv \tau^2 \kappa$, and the terms involving \mathbf{X} drop out, with estimation of β absorbed into \mathbf{g} , thereby avoiding non-identifiability.

One can simplify the expected MSE analytically, making use of the fact that $\mathbf{S} \mathbf{K} \mathbf{X} = \mathbf{X}$ (because $\mathbf{Q} \mathbf{X} = \mathbf{0}$ when the columns of \mathbf{X} are in the space spanned by the eigenvectors of \mathbf{Q}):

$$\begin{aligned} \text{MSE}(\lambda) &= E_g E_Y ((\hat{\mathbf{g}} - \mathbf{g})^\top (\hat{\mathbf{g}} - \mathbf{g})) \\ &= E_g E_Y (\mathbf{Y}^\top \mathbf{S}^\top \mathbf{S} \mathbf{Y} - 2 \mathbf{g}^\top \mathbf{S} \mathbf{Y} + \mathbf{g}^\top \mathbf{g}) \\ &= E_g (\tau^2 \mathbf{S}^\top \mathbf{S} - 2 \mathbf{g}^\top \mathbf{S} \mathbf{K} \mathbf{g} + \mathbf{g}^\top \mathbf{g} + \mathbf{g}^\top \mathbf{K}^\top \mathbf{S}^\top \mathbf{S} \mathbf{K} \mathbf{g}) \\ &= \tau^2 \text{tr}(\mathbf{S}^\top \mathbf{S}) - 2 \text{tr}(\mathbf{S} \mathbf{K} \mathbf{C}) + \text{tr}(\mathbf{C}) + \text{tr}(\mathbf{K}^\top \mathbf{S}^\top \mathbf{S} \mathbf{K} \mathbf{C}) \\ &\quad - 2(\mathbf{X} \beta)^\top \mathbf{S} \mathbf{K} \mathbf{X} \beta + (\mathbf{X} \beta)^\top \mathbf{X} \beta + (\mathbf{X} \beta)^\top \mathbf{K}^\top \mathbf{S}^\top \mathbf{S} \mathbf{K} \mathbf{X} \beta \\ &= \tau^2 \text{tr}(\mathbf{S}^\top \mathbf{S}) - 2 \text{tr}(\mathbf{S} \mathbf{K} \mathbf{C}) + \text{tr}(\mathbf{C}) + \text{tr}(\mathbf{K}^\top \mathbf{S}^\top \mathbf{S} \mathbf{K} \mathbf{C}). \end{aligned} \tag{6}$$

Note that the MSE can be expressed in terms of a single parameter to estimate, λ , plus the known parameters of the data-generating model. To simplify the presentation of results, I consider an oracle result based on optimizing the MSE (6) across all possible values of λ for each of the ICAR and TPS-MRF models, thereby using the best overall λ for a given generative setting.

Analogous calculations for MSE can be done when fitting the model using the GP approach (therby fitting with the same model used to generate observations), giving (6) but with \mathbf{S} involving \mathbf{R}_θ^{-1} in place of \mathbf{Q} and $\lambda \equiv \frac{\tau^2}{\sigma^2}$. Since the GP is the generating model, it provides a baseline for comparison with the MRF models, so I also assume that $\lambda = \frac{\tau^2}{\sigma^2}$ is known, as are $\boldsymbol{\theta} = \{\rho, \nu\}$ and $\beta = \mathbf{0}$.

3.2.1 Point observations

Fig. 5a shows the ratio of the MSE using the ICAR model to that using the TPS model as well as the ratio of the MSE using the true GP model to that using the TPS model. For $n = 100$, representing fairly sparse data, we see that with uniformly distributed locations, in general the TPS model either matches the MSE of the ICAR or improves upon it. The TPS model strongly outperforms the ICAR when $\nu = 2$, the range is moderate to large, and the noise variance is not too large. The one case in which the ICAR beats the TPS-MRF is when $\rho = 0.08$, particularly for smaller values of τ^2 and $\nu = 2$, but note that in this case the absolute MSE (red lines) is large in for both models. For $\nu = 0.5$, which produces locally heterogeneous surfaces for which we would expect the ICAR model to perform well, we see that the two MRF models perform fairly similarly. With clustered data, the results are qualitatively similar, but we see that the TPS-MRF has less of an advantage over the ICAR model. This is likely a result of the TPS-MRF approximating a spline, with no constraint on the spline fit in the large gaps with no observations. Unlike a Gaussian process, a spline does not revert to an overall estimate of the mean when it is far from any observations, which may lead to poor interpolation. Comparing the TPS-MRF to the GP, the GP generally performs better, which is not surprising given that it is the generative model and is fit based on the true hyperparameter values, but particularly for uniformly-sampled locations, the

TPS-MRF is competitive.

Results are generally similar for $n = 1000$ (not shown) except that with more data, the value of ρ at which the TPS-MRF starts to outperform the ICAR model is 0.08, compared to 0.32 when $n = 100$. The effect of clustering is also more detrimental to the TPS-MRF, particularly at $\rho = 0.08$. In addition the TPS model shows more advantage over the ICAR at the larger τ^2 values for $\nu = 2$.

I interpret these results as follows. For small values of the range, the oracle fit is a constant surface, and both MRF models do this and perform similarly, with MSE approximately equal to $E_g \mathbf{g}^\top \mathbf{g}$, the MSE of $\hat{\mathbf{g}} = \mathbf{0}$; see the example fits in Fig. 6, first column. Neither model has any predictive ability, and MSE is essentially squared bias. For slightly larger range values, the ICAR model does modestly better (example fits in Fig. 6, second column), with the TPS-MRF oversmoothing (and experiencing some boundary issues) and the more local behavior and mean reversion of the ICAR model being beneficial. Then, as the range increases the TPS improves upon the ICAR, most notably when the true surface is smooth ($\nu = 2$) (Fig. 6, third column), with the ICAR showing the bulls-eyes seen in Fig. 1. If the noise variance is also very large, then the advantage of the TPS when $\nu = 2$ moderates, with the TPS oversmoothing (Fig. 6, fourth column). When $\nu = 0.5$, performance of the ICAR and TPS-MRF are comparable because the TPS oversmooths, but follows the larger-scale patterns better than the ICAR, while the ICAR more closely follows the local variability in the vicinity of the observations (not shown).

Fig. 5b shows results when there is one observation per grid cell, a pure smoothing problem without any interpolation. Except for the smallest value of ρ and largest τ^2 , the TPS does as well as the ICAR for $\nu = 0.5$ and outperforms the ICAR model for $\nu = 2$, while matching the performance of the GP.

3.2.2 Areal observations

Finally, for areal data, the MSE for the TPS-MRF model is similar to the ICAR and GP, with the TPS-MRF outperforming the ICAR model when τ^2 is smaller and for smoother surfaces (Fig. 5c),

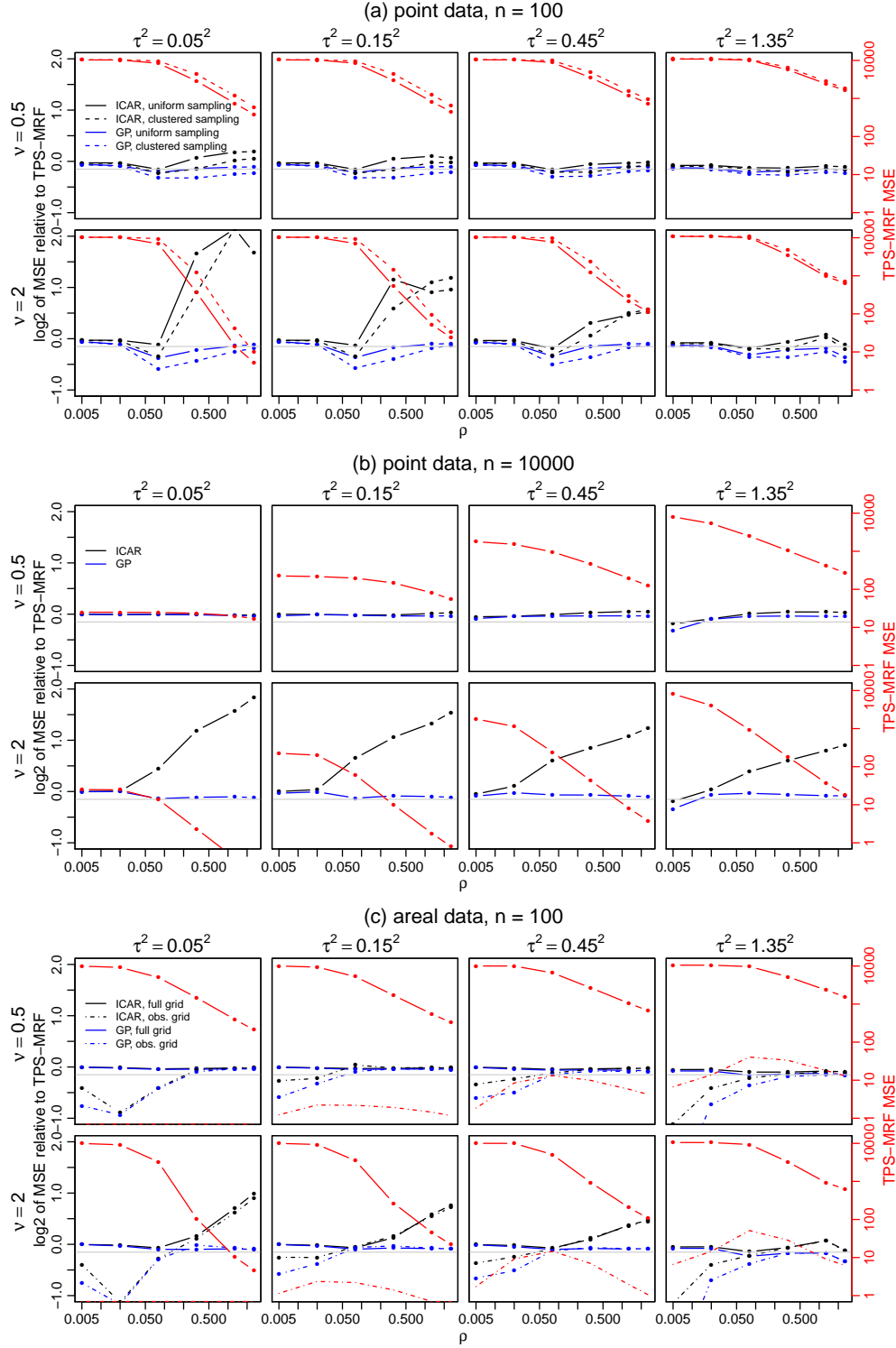


Figure 5: Predictive performance for (a) $n = 100$ point observations, (b) $n = 10000$ point observations, and (c) $n = 100$ areal observations; all are based on oracle values for λ for the MRF models and the true parameter values for the GP estimation. Plots show the log (base 2) of the ratio of MSE for the ICAR (black) and GP (blue) models relative to the TPS-MRF and absolute MSE for the TPS-MRF for reference (in red, with axis labels on the right side). In each subplot, ν varies with the row and τ^2 with the column. The horizontal grey line corresponds to the MSE being 90% of the MSE of the TPS-MRF as an informal cutoff below which other models perform substantially better.

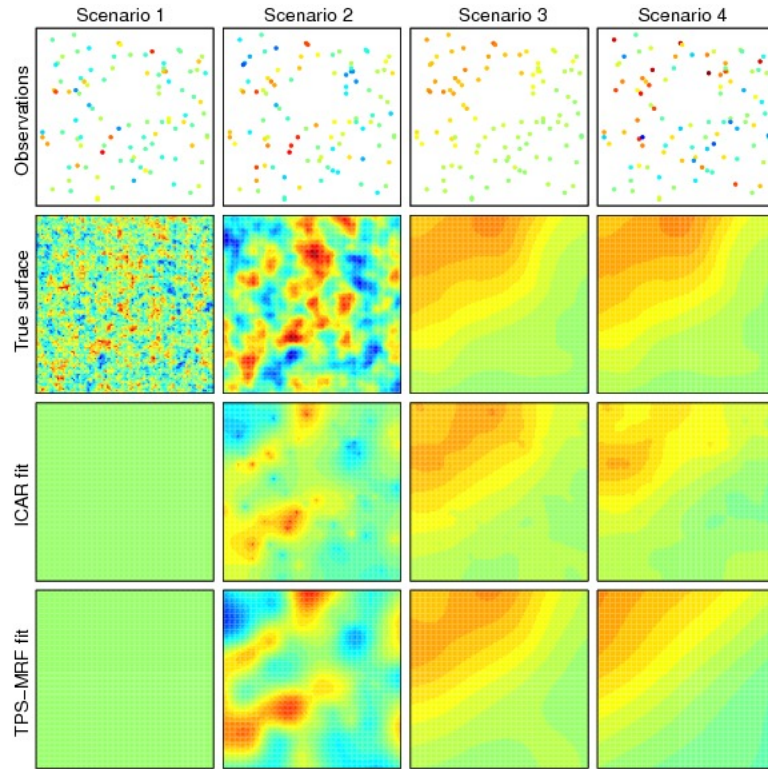


Figure 6: Example ICAR and TPS-MRF fits for four generative scenarios with uniform sampling of locations: (1) $\nu = 0.5$, $\rho = 0.02$, $\tau^2 = 0.15^2$, (2) $\nu = 2$, $\rho = 0.08$, $\tau^2 = 0.15^2$, (3) $\nu = 2$, $\rho = 1.28$, $\tau^2 = 0.15^2$, (4) $\nu = 2$, $\rho = 1.28$, $\tau^2 = 1.35^2$.

similar to the results for point data. When considering the MSE of predictions of the surface integrated to the resolution of the coarse grid at which the observations are available (dotted lines), we see that for smaller values of the range, the TPS-MRF performs worse than the ICAR and the GP. However, note that in these cases there is little spatial signal in the areal data (i.e., in the data at the observed resolution).

3.3 Simulations

The oracle results presume good choices of the penalty parameter for each representation and do not assess the effect of hyperparameter estimation on performance. Here I present basic simulation results under a set of scenarios informed by the analytic results. I consider five scenarios: $n = 100$ and $n = 1000$, each with points sampled uniformly from the domain and sampled in a clustered fashion, and $n = 100$ areal observations on the coarse 10 by 10 grid. For each scenario I carry out a full factorial design with respect to the data-generating parameters: $\nu \in \{0.5, 2\}$, $\rho \in \{0.04, 0.16, 0.64\}$, and $\tau^2 \in \{0.15^2, 0.45^2, 1.35^2\}$. Otherwise the generative model is as in Section 3.2. For the areal scenario, I do not fit the GP model because of computational constraints. I use 100 simulations; the resulting Monte Carlo standard errors are small relative to the reported MSE results.

The unknowns are \mathbf{g} , τ^2 , and $\lambda = \kappa\tau^2$ (for the MRF) or $\lambda = \tau^2/\sigma^2$ (for a GP model) and β (for the GP model). Integrating over \mathbf{g} to obtain a marginal likelihood in terms of τ^2 and λ and profiling over τ^2 (and β for the GP) gives a likelihood that can be numerically maximized with respect to λ . One can then use (5) to estimate \mathbf{g} (with the analogous quantity for the GP case). In fitting the GP model, I use the true ν .

3.3.1 Point observations

The TPS-MRF often outperforms both the ICAR and GP in fitting point data, in particular for the larger two values of ρ (Fig. 7a,b). The primary exception is for $\rho = 0.04$ with $n = 1000$. There are only limited differences when comparing uniform sampling to clustered sampling of locations.

In the simulations the TPS-MRF outperforms the GP in some scenarios, in contrast to the oracle assessment (Section 3.2) that used the true GP parameter values. The results suggest that while the use of the TPS model can hurt performance somewhat when the true surface is quite wiggly, it can give substantial improvements when the surface is smoother. Results are shown for MSE within the convex hull of the observations; if one considers the full spatial domain, the patterns are similar, but the relative performance of the TPS-MRF is not as good, indicating boundary issues for the TPS-MRF.

In terms of uncertainty characterization, both the ICAR and TPS-MRF give the nominal 95% coverage for prediction intervals. However, coverage for the true function, g , is quite low in some situations. In particular, for $n = 100$ both the ICAR and TPS-MRF have coverage below 70% for $\rho = 0.04$. For $\rho = 0.16$, the TPS-MRF shows low coverage for all values of ν and τ^2 and for $\rho = 0.64$ when $\nu = 0.5$. The good predictive coverage and poor function coverage of the TPS-MRF occurs because the model assumes a smooth underlying surface and attributes some of the true variability in the function surface to the error component, thereby giving standard errors for the function values that are too small. Banerjee et al. (2010) noted a similar problem for reduced rank kriging, with an inflation in the estimated nugget variance. This suggests that one interpret the TPS-MRF uncertainty as relating to the larger-scale spatial variability and accept that one is not able to characterize uncertainty about finer-scale variation.

3.3.2 Areal observations

Fig. 7c shows results for areal data. The TPS-MRF does better when $\nu = 2$ for large ρ , while the ICAR performs better for $\rho = 0.16$, particularly when the noise level is low. For the degree of spatial structure specified with $\rho = 0.16$, at the aggregated resolution of the observations (results shown in dashed lines) the true surface is fairly heterogeneous between neighboring areas and the ICAR follows the data closely, while the TPS-MRF oversmooths.

Coverage of the TPS-MRF prediction intervals is at the nominal 95% level, but for large values of ρ and smaller values of τ^2 , the ICAR model often follows the data exactly, estimating $\tau^2 = 0$

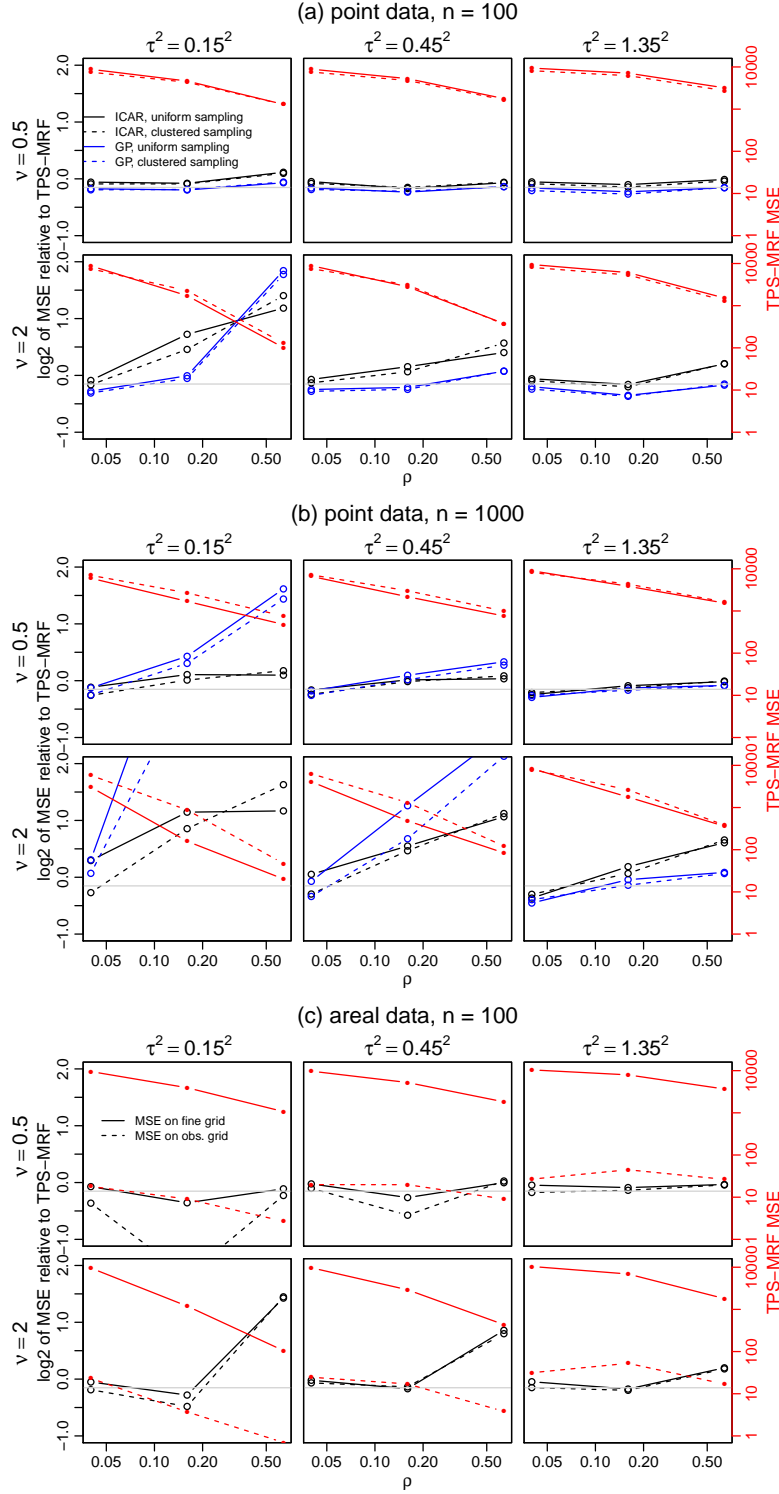


Figure 7: Predictive performance based on simulation for (a) $n = 100$ point observations, (b) $n = 1000$ point observations, and (c) $n = 100$ areal observations. Predictive accuracy is shown as the log (base 2) of the average over 100 simulations of the ratio of MSE for either ICAR (black) or (for point observations only) GP (blue) models relative to TPS-MRF, with the average absolute MSE for the TPS-MRF shown in red for reference (with axis labels on the right side). For point observations, MSE is computed for locations within the convex hull of the observations in a given simulation. The horizontal grey line corresponds to the MSE being 90% of the MSE of the TPS-MRF as an informal cutoff below which other models perform substantially better.

and giving very low coverage (not shown). Considering coverage for the true \mathbf{g} , results are similar to the situation for point observations for both models, with coverage below 50% when $\rho = 0.04$. However, note that in this scenario there is little spatial signal in the aggregated observations. The TPS-MRF shows low coverage generally, except when $\nu = 2$, $\rho = 0.64$, and $\tau^2 \leq 0.45^2$, i.e., coverage is good only when it is easy to follow the observations. The coverage for the ICAR model is generally better, but conservative when both $\rho \geq 0.16$ and $\tau^2 \leq 0.45^2$. However, the ICAR coverage is generally below 50% in the presence of little noise, again because it estimates that $\tau^2 = 0$.

4 Computational considerations

4.1 Normal data

Consider an MRF model for \mathbf{g} , $\mathbf{g} \sim \mathcal{N}_m(\mathbf{0}, (\kappa\mathbf{Q})^-)$, where the zero mean is justified by (5). If the observations are normally distributed, $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta} + \mathbf{K}\mathbf{g}, \tau^2\mathbf{I})$, we have conjugacy and can integrate \mathbf{g} out of the model to obtain a marginal likelihood with which to do maximization or MCMC on the hyperparameter space. Note that one could also use the INLA methodology to quickly approximate the posterior without MCMC (Rue et al., 2009), but here I explore the computations needed for maximum likelihood and for 'exact' inference via MCMC.

The marginal precision for \mathbf{Y} can be expressed, based on the Sherman-Morrison-Woodbury formula, as

$$\boldsymbol{\Sigma}_\lambda^{-1} = \frac{1}{\tau^2} \left(\mathbf{I} - \mathbf{K}(\lambda\mathbf{Q} + \mathbf{K}^\top\mathbf{K})^{-1}\mathbf{K}^\top \right).$$

For maximization, we can express the maxima for $\boldsymbol{\beta}$ and τ^2 as functions of $\boldsymbol{\Sigma}_\lambda^{-1}$ and therefore of λ ,

$$\begin{aligned} \hat{\boldsymbol{\beta}}_\lambda &= (\mathbf{X}^\top \boldsymbol{\Sigma}_\lambda^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}_\lambda^{-1} \mathbf{Y} \\ \hat{\tau}_\lambda^2 &= \frac{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_\lambda)^\top \boldsymbol{\Sigma}_\lambda^{-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_\lambda)}{n - c}, \end{aligned}$$

where c is the number of zero eigenvalues of \mathbf{Q} . Both of these quantities can be calculated efficiently based on a sparse Cholesky decomposition of $\lambda\mathbf{Q} + \mathbf{K}^\top \mathbf{K}$ in the expression for Σ_λ^{-1} , because $\mathbf{K}^\top \mathbf{K}$ will generally be sparse. If all the data are point locations, $\mathbf{K}^\top \mathbf{K}$ is diagonal, with the diagonal entries counting the number of observations falling in each grid cell. If all the data are areal observations, only off-diagonal elements of $\mathbf{K}^\top \mathbf{K}$ corresponding to pairs of grid cells that are overlapped by a common areal observation are non-zero.

The marginal profile likelihood as a function of λ alone is proportional to

$$\frac{\lambda^{(m-c)/2}}{(\hat{\tau}_\lambda^2)^{(n-c)/2} |\lambda\mathbf{Q} + \mathbf{K}^\top \mathbf{K}|^{1/2}}$$

where the determinant can be calculated efficiently based on the already-computed sparse Cholesky decomposition.

MCMC calculations rely on similar quantities that can be computed efficiently. To draw from the posterior of \mathbf{g} off-line, given κ , τ^2 , and β , we have

$$\mathbf{g} \sim \mathcal{N}((\mathbf{K}^\top \mathbf{K} + \kappa\tau^2\mathbf{Q})^{-1} \mathbf{K}^\top (\mathbf{Y} - \mathbf{X}\beta), \tau^2(\mathbf{K}^\top \mathbf{K} + \lambda\mathbf{Q})^{-1}),$$

which can be done efficiently based on the same Cholesky decomposition as above.

The computational limitation in the proposed MRF approach is the ability to work with a sparse matrix, $\mathbf{K}^\top \mathbf{K} + \lambda\mathbf{Q}$, whose size scales with the number of grid cells. Thus the approach is limited only by the resolution at which one wishes to do prediction rather than by the sample size, which is similar to the computational constraint in reduced rank kriging, where the computational cost scales with the number of knots. For small n and large m , it may be more computationally efficient to represent the data covariance as

$$\Sigma = \tau^2 \left(\mathbf{I} + \frac{1}{\lambda} \mathbf{K} \mathbf{Q}^{-1} \mathbf{K}^\top \right)$$

and precompute the n by n matrix $\mathbf{K} \mathbf{Q}^{-1} \mathbf{K}^\top$, while in the iterations of an MCMC or optimization

computing the Cholesky of the dense matrix Σ . Note that here we need to use the generalized inverse, thereby assigning zero variance to the eigenvectors of \mathbf{Q} corresponding to zero eigenvalues and hence necessitating inclusion of the relevant terms in the mean of \mathbf{g} .

4.2 Non-normal data

The general model (1-2) is a GLMM, where \mathbf{K} and \mathbf{g} play the roles of the random effects design matrix and random effects, respectively. In this case the covariance of the random effects has spatial structure and is specified in terms of a precision matrix. Both likelihood-based and Bayesian inference for GLMMs is computationally challenging because inference involves a high-dimensional integral with respect to the random effects that cannot be expressed in closed form.

4.2.1 PQL

The PQL approach (Breslow and Clayton, 1993; Wolfinger and O’Connell, 1993) maximizes a GLMM log-likelihood using a Laplace approximation to the integral over the random effects. The PQL method iterates between optimizing with respect to the hyperparameters, based on the approximation to the marginal likelihood, and using an iterated weighted least squares approach to optimize with respect to both the fixed and random effects.

While PQL is implemented in R in the `glmmPQL()` function in the MASS package, `glmmPQL()` will not fit the models described here because the `corStruct` classes for the lme package do not include MRF specifications in two dimensions. Furthermore, the function is not set up to work efficiently with sparse precision matrices. PQL is also implemented in R in the new `glmmGS` package, which implements the Gauss-Seidel optimization approach of Guha et al. (2009). However, while `glmmGS` allows arbitrary precision matrices that can be specified and manipulated as sparse matrices, at present \mathbf{K} cannot be specified to be a sparse matrix, so the implementation is not computationally feasible for large numbers of grid cells. Thus, custom code is currently required to implement the models discussed here.

4.2.2 Integrated Nested Laplace Approximation (INLA)

(Rue et al., 2009) present an approach to fitting GLMMs based on nested Laplace approximations involving both the hyperparameters and the latent process values. The result is estimation of the marginal posterior densities of the random effects and the hyperparameters. Note that this accounts for uncertainty in hyperparameters, in contrast to the maximization done with the PQL approach. The INLA R package (www.r-inla.org) can make use of sparsity in both \mathbf{Q} and \mathbf{K} and is therefore very computationally efficient. For analyses using likelihood and prior models that are implemented in INLA and for which one needs only marginal posteriors (or posteriors of linear combinations), INLA is a promising option.

4.2.3 MCMC

While MCMC is a standard approach for Bayesian GLMMs, convergence and mixing are often troublesome (Christensen and Waagepetersen, 2002; Christensen et al., 2006), because of the high-dimensionality of the random effects, the dependence between random effects (particularly when these represent spatial or temporal structure), and cross-level dependence between random effects and their hyperparameters ($\{\tau^2, \kappa\}$ in this work) (Rue and Held, 2005; Rue et al., 2009). The sparse matrix calculations possible with MRF models improve computational efficiency but do not directly address convergence and mixing issues.

Gamerman (1997) describes the use of a weighted least squares proposal for the fixed and random effects, also suggested in Rue and Held (2005, pp. 167-169) to deal with the dependence amongst the random effects, and Rue and Held (2005) suggest combining this with a joint update of the hyperparameter and the random effects to address the cross-level dependence. For high-dimensional random effects vectors such as proposed here, using subblocks of \mathbf{g} may also be helpful. To simulate draws of \mathbf{g} from the posterior, one might also use PQL or INLA to estimate the hyperparameters and then conditionally draw samples of the fixed and random effects using MCMC.

One disadvantage of the TPS-MRF compared to the ICAR may be that the additional smooth-

ness makes it more difficult to accept MCMC proposals for \mathbf{g} because the value of the process in one grid cell is more strongly constrained by the values in the other grid cells. The Brownian motion-like behavior in the ICAR (Besag and Mondal, 2005) may help to decouple values of \mathbf{g} for different grid cells, similar to the strategy of adding a small amount of white noise to the latent process (e.g., Wikle, 2002), which Paciorek (2007) showed could improved MCMC mixing in a GP context.

5 Examples

5.1 Point-level pollution modeling

This example is based on the work of Paciorek et al. (2009), who modeled spatio-temporal variation in air pollution for the purpose of predicting concentrations for use as the exposure values in a health analysis. The data are average fine particulate matter over 2001-2002 at 339 monitoring stations in the northeast U.S., from the US EPA’s Air Quality System database. For this analysis, I averaged the 24 monthly average values and included only locations with at least 22 months of data. I compare the use of the MRF approach for modeling point-level data (based on a 4 km resolution grid) with an additive model built on a reduced rank thin plate spline (using the `gam()` function from the `mgcv` package in R) and to universal kriging with both an exponential covariance and a Matérn covariance with $\nu = 2$. For the MRF models, I fitted a model of the form (1-2) with a linear link and normal likelihood (with independent, homoscedastic errors) using the computational approach described in Section 4.1. I used the same set of covariates (log of the distance to nearest road in two road size classes, percent urban land use in a local buffer, log of elevation, and log of the estimated fine PM emissions within a 10 km buffer) as in Paciorek et al. (2009). I included the covariates as linear terms for simplicity and because Paciorek et al. (2009) found only a minor improvement when considering an additive nonparametric structure for the covariates. I used ten-fold cross-validation to assess hold-out error for 219 stations, with 120 stations in boundary states forming a spatial buffer and always used in the training set, to compare

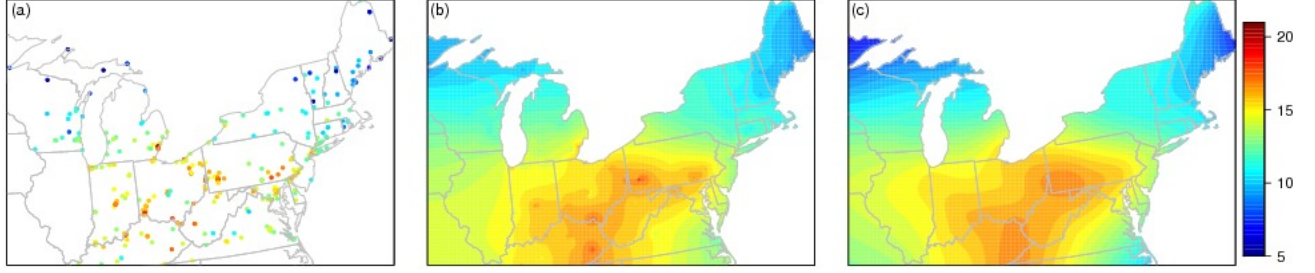


Figure 8: PM observations (a) and fitted residual spatial surfaces (not including the effect of covariates) for the ICAR (b) and TPS-MRF (c) models.

the models.

The models all gave very similar prediction results, with mean squared prediction errors of 1.30 for the `gam()` function, 1.34 and 1.29 for universal kriging with the Matérn ($\nu = 2$) and exponential covariances respectively, 1.32 for the TPS-MRF and 1.26 for the ICAR. These correspond to cross-validation R^2 values of about 0.8 (0.787-0.811). The slight apparent advantage for kriging with an exponential covariance and for the ICAR model suggests the presence of fine-scale variability that the other models smooth over. Fig. 8 shows the estimates of residual spatial variability, not including the effect of the spatially-varying covariates, illustrating that the ICAR model fits local effects around the observations but also seems to generally follow the large-scale pattern seen in the TPS-MRF. All the models showed good hold-out prediction coverage, with the MRF models having larger standard errors and therefore some overcoverage. Interestingly, in light of the scale dependence of the oracle and simulation results, in this setting in which it appears there is variation at multiple scales, the ICAR appears to slightly outperform the TPS-MRF, which may also relate to the fact that the observations are clustered (in metropolitan areas), a setting that the oracle and simulation results suggest works in favor of the ICAR model.

Paciorek and Liu (2012) considered more complicated models of fine PM over a variety of spatial and temporal domains using the MRF approach outlined in this paper, including combining point-level monitoring data with areal data from remote sensing and fitting a spatio-temporal extension of the model proposed here.

5.2 Area-level disease mapping

This example is based on the work of Krieger et al. (2006) and Hund et al. (2012), who analyzed variation in breast cancer incidence in Los Angeles County, California, US. Their analyses focused on the relationship between census tract-level poverty and breast cancer incidence and changes over time in that relationship, while accounting for residual spatial variation. Here I modeled breast cancer incidence for data from 1998-2002 for white non-Hispanic women.

The data are counts of cancer incidence in the 5-year period. Following Krieger et al. (2006) and Hund et al. (2012), I fit Poisson models with a log link, using the log of the expected counts as an offset term,

$$\log \mu_i = \log E_i + \beta_0 + \boldsymbol{\beta}^\top \mathbf{pov}_i + \mathbf{K}_i^\top \mathbf{g},$$

where μ_i and E_i are the Poisson mean and the expected number of cases in the i th census tract (CT). The expected numbers were calculated based on internal age standardization, described in Hund et al. (2012) and based on CT population (multiplied by five, which assumes constant population over the 5-year period) from the 2000 US Census. CT poverty was a five-level categorical variable with \mathbf{pov}_i being a vector of four indicator variables determining the poverty category of the i th CT. The categorical poverty variable is defined as follows: (1) <5% of residents living below the poverty line and more than 10% of households having high income (at least four times the US median household income), (2) <5% of residents living below the poverty line and less than 10% of households having high income, (3) 5.0-9.99% of residents living below the poverty line, (4) 10.0-19.99% of residents living below the poverty line, and (5) at least 20% of residents living below the poverty line, which was used as the baseline category.

I used the INLA package in R to fit the ICAR and TPS-MRF models defined on a fine grid with 99×101 cells of size 1.25^2 km^2 and compared the results to a standard ICAR model based on the neighborhood structure of the irregular census tracts, in all cases using the default hyperparameter priors specified in INLA. The estimated coefficients for the categorical poverty variable were very similar to those in Hund et al. (2012), with lower poverty CTs showing higher breast cancer inci-

dence. Fig. 9 shows the log of the raw incidence rate ratio (observed counts divided by expected number), compared to the estimated log-incidence rate ratios, \hat{g} , for the three models. Values of zero indicate no departure from the expected number based on the population and age distribution in the CT. The census-tract-based ICAR and fine grid-based ICAR models show much more spatial variability, while the TPS-MRF smooths quite a bit. DIC and the summed log conditional predictive ordinate (CPO) values, $\sum_i \log P(Y_i|Y_{-i})$ (provided by the INLA package) suggest that the census-tract based ICAR (DIC of 8438, logCPO of -4227) outperforms the grid-based ICAR (8451, -4233) and TPS-MRF models (8476, -4241). The advantage of the ICAR over the TPS-MRF is consistent with the simulation results for areal data in which the ICAR performed better for the moderate value of ρ (Fig. 7c).

6 Discussion

I have presented a straightforward modeling approach for both areal and point data that relates observations to an underlying smooth spatial surface, represented as an MRF. One important result is that the analytic comparison of various MRF structures indicates that higher-order neighborhood structures that do not have the weighting structure of the TPS-MRF do not produce smoother processes than the standard ICAR model. This suggests that such models are not appealing for spatial modeling. Given these results it would be useful to investigate the smoothing properties of models that empirically choose neighborhood structures (White and Ghosh, 2009; Zhu et al., 2010).

Based on the analytic and simulation assessment of predictive performance, the TPS-MRF outperforms the ICAR model in many scenarios, in particular with smoother surfaces (both in terms of the spatial range and differentiability), as would be expected given that the TPS-MRF approximates a thin plate spline. In both examples, however, the ICAR appeared to outperform the TPS-MRF, perhaps because the true surface was relatively wiggly in those contexts. One open question is how the models would perform in a setting with variation at multiple resolutions; these

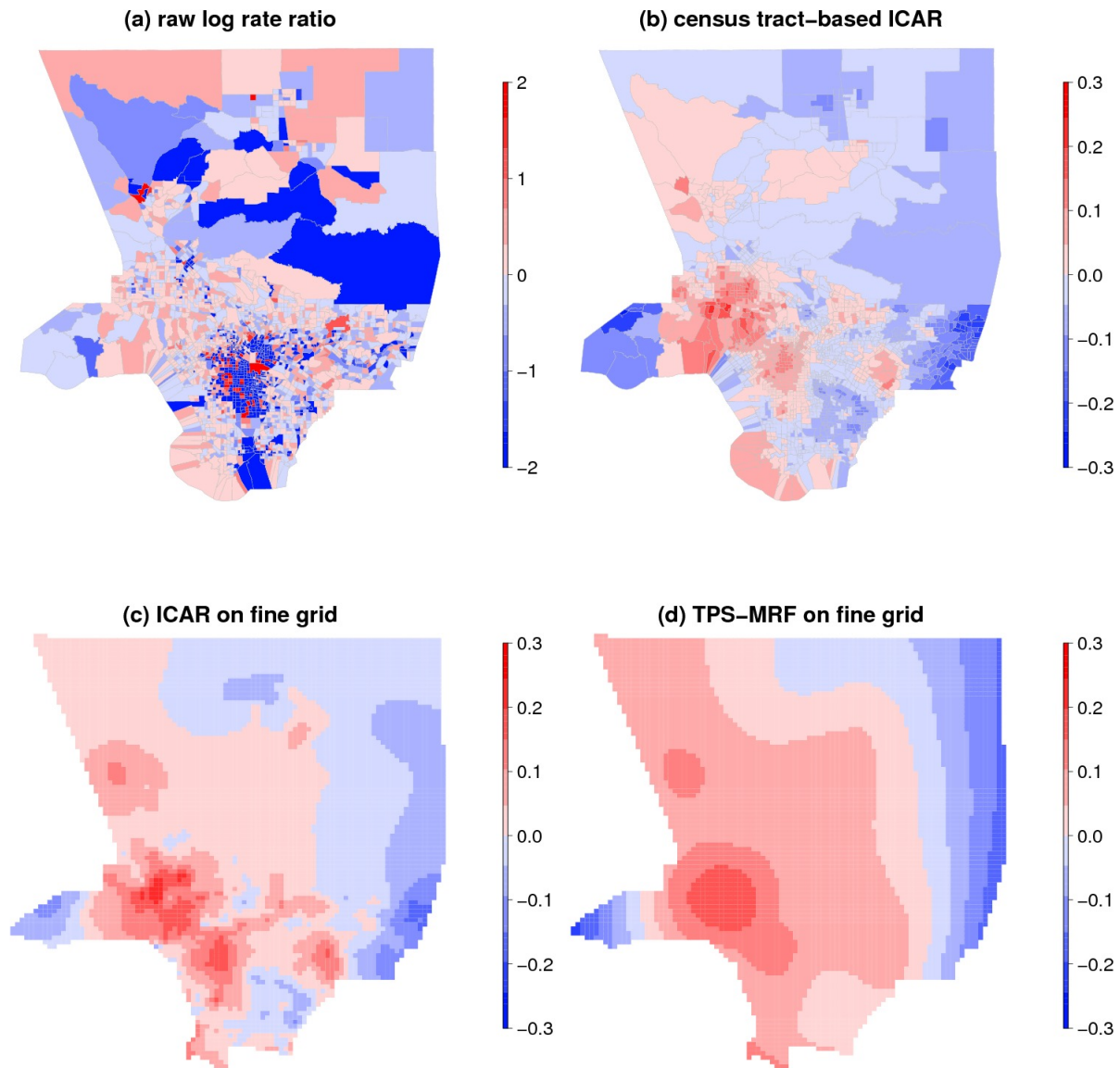


Figure 9: (a) Raw log breast cancer incidence rate ratios by census tract in Los Angeles County and estimated log incidence rate ratios for (b) a standard ICAR model based on the census tract neighborhood structure, (c) the ICAR model on a fine grid, and (d) the TPS-MRF model on a fine grid. Note that the scale in (a) is different than the other panels and that values with magnitude greater than two are censored. Also note that the census tracts extend into the Pacific Ocean in the southwestern portion of panels (a) and (b), while panels (c) and (d) show only the land-based grid cells.

results suggest that the TPS-MRF better represents the large scale while the ICAR model better captures fine-scale variation. How these balance in terms of overall error would likely depend on the relative magnitude of the variation at the different scales.

Spline models can do poorly in situations with large spatial gaps (where large is relative to the spatial range of dependence in the process being modeled) and on the boundary of the domain, as the estimation of basis coefficients is poorly constrained by the data and influenced by data at the extremes of the support of the basis functions. The TPS-MRF model, by virtue of approximating a thin plate spline, can have this unappealing behavior. In contrast a GP model, being a stationary model, gives predictions that revert to the overall mean for prediction locations far (relative to the estimated spatial range) from observations. In many large datasets, for which the computational efficiency of MRF models is appealing, including deterministic model output and remote sensing observations, gaps are not present or tend to be small, so the issue of extrapolation into large gaps may not be a concern. Furthermore, my ad hoc experience suggests that gaps are less of a problem in two dimensions than in one dimension, although the TPS-MRF can have problems at the boundaries of the domain.

An appealing alternative to the MRF models presented here is the MRF construction of Lindgren et al. (2011), which approximates a Gaussian process with Matérn covariance for integer values of the Matérn smoothness (differentiability) parameter. I expect that much future work with MRFs will involve this construction because of the added flexibility of a representation that includes the GP range parameter. However, I note that the simulations suggest that the TPS-MRF in some cases outperforms an exact GP representation, and in the LA example, I found that the ICAR models outperformed the Lindgren et al. (2011) model (results not shown). Lindgren et al. (2011) focus on a triangulation rather than a rectangular grid, which has computational advantages when dealing with irregular domains.

I have highlighted the advantages of using a smooth underlying surface for areal data. These include the ability to deal with data aggregation in a consistent manner and with spatial misalignment. Furthermore, in many situations, the area boundaries are essentially arbitrary relative to

the process being measured, so the resulting neighborhood structure is arbitrary as well. Rather it is appealing to imagine an smooth underlying surface, with the areal units merely a measurement artefact that is represented in the measurement model through the mapping matrix, \mathbf{K} . In some cases, administrative units might actually have a direct effect on the outcome, in which case more traditional MRF models based on a single random effect per area and standard neighborhood structures may be more appealing, although independent random effects may be appealing in some cases.

Spatio-temporal modeling situations are of course very common. Paciorek and Liu (2012) describes a spatio-temporal extension of the spatial models described here that allows for autoregressive structure in time. In contrast to many spatio-temporal models, the approach has a spatial mean that is shared across time points, which Stein and Fang (1997) emphasize is important for allowing one to properly characterize uncertainty when aggregating over time periods.

Acknowledgments

The author thanks Jeff Yanosky for the particulate matter dataset, Jarvis Chen, Lauren Hund, and Brent Coull for access to the LA breast cancer dataset, Steve Melly for GIS processing for the LA example, and Finn Lindgren for helpful discussions. This work was funded by NCI P01 Grant CA134294-01.

References

- Banerjee, S., Finley, A., Waldmann, P., and Ericsson, T. 2010. Hierarchical spatial process models for multiple traits in large genetic trials. *Journal of the American Statistical Association* 105: 506–521.
- Banerjee, S., Gelfand, A., Finley, A., and Sang, H. 2008. Gaussian predictive process models for large spatial data sets. *Journal of the Royal Statistical Society: Series B* 70: 825–848.

- Banerjee, S., Gelfand, A., and Sirmans, C. 2003. Directional rates of change under spatial process models. *Journal of the American Statistical Association* 98: 946–954.
- Besag, J. and Mondal, D. 2005. First-order intrinsic autoregressions and the de Wijs process. *Biometrika* 92: 909–920.
- Breslow, N. E. and Clayton, D. G. 1993. Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association* 88: 9–25.
- Christensen, O., Roberts, G., and Sköld, M. 2006. Robust Markov chain Monte Carlo methods for spatial generalized linear mixed models. *Journal of Computational and Graphical Statistics* 15: 1–17.
- Christensen, O. F. and Waagepetersen, R. 2002. Bayesian prediction of spatial count data using generalized linear mixed models. *Biometrics* 58: 280–286.
- Diggle, P., Menezes, R., and Su, T.-L. 2010. Geostatistical inference under preferential sampling. *Journal of the Royal Statistical Society, Series C* 59: 191–232.
- Fuentes, M. and Raftery, A. 2005. Model evaluation and spatial interpolation by Bayesian combination of observations with outputs from numerical models. *Biometrics* 61: 36–45.
- Furrer, R., Genton, M. G., and Nychka, D. 2006. Covariance tapering for interpolation of large spatial datasets. *Journal of Computational and Graphical Statistics* 15: 502–523.
- Gamerman, D. 1997. Sampling from the posterior distribution in generalized linear mixed models. *Statistics and Computing* 7: 57–68.
- Guha, S., Ryan, L., and Morara, M. 2009. Gauss–Seidel estimation of generalized linear mixed models with application to Poisson modeling of spatially varying disease rates. *Journal of Computational and Graphical Statistics* 18: 818–837.
- Hrafnkelsson, B. and Cressie, N. 2003. Hierarchical modeling of count data with application to nuclear fall-out. *Environmental and Ecological Statistics* 10: 179–200.

- Hund, L., Chen, J., Krieger, N., and Coull, B. 2012. A geostatistical approach to large-scale disease mapping with temporal misalignment. *Biometrics* in press.
- Kamman, E. and Wand, M. 2003. Geoadditive models. *Applied Statistics* 52: 1–18.
- Kaufman, C., Schervish, M., and Nychka, D. 2008. Covariance tapering for likelihood-based estimation in large spatial datasets. *Journal of the American Statistical Association* 103: 1556–1569.
- Kelsall, J. and Wakefield, J. 2002. Modeling spatial variation in disease risk. *Journal of the American Statistical Association* 97: 692–701.
- Krieger, N., Chen, J., Waterman, P., Rehkopf, D., Yin, R., and Coull, B. 2006. Race/ethnicity and changing us socioeconomic gradients in breast cancer incidence: California and Massachusetts, 1978–2002 (United States). *Cancer Causes and Control* 17: 217–226.
- Lindgren, F., Rue, H., and Lindström, J. 2011. An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society, Series B* 73: 423–498.
- Mugglin, A., Carlin, B., and Gelfand, A. 2000. Fully model-based approaches for spatially misaligned data. *Journal of the American Statistical Association* 95: 877–887.
- Paciorek, C. 2007. Bayesian smoothing with Gaussian processes using Fourier basis functions in the spectralGP package. *Journal of Statistical Software* 19: 2.
- 2012. Combining spatial information sources while accounting for systematic errors in proxies. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 61: 429–451.
- Paciorek, C. and Liu, Y. 2012. Assessment and statistical modeling of the relationship between remotely-sensed aerosol optical depth and PM_{2.5}. *Health Effects Institute Research Report 167 (peer-reviewed)*.

- Paciorek, C., Yanosky, J., Puett, R., Laden, F., and Suh, H. 2009. Practical large-scale spatio-temporal modeling of particulate matter concentrations. *Annals of Applied Statistics* 3: 369–396.
- Pettitt, A., Weir, I., and Hart, A. 2002. A conditional autoregressive Gaussian process for irregularly spaced multivariate data with application to modelling large sets of binary data. *Statistics and Computing* 12: 353–367.
- Rue, H. and Held, L. 2005. *Gaussian Markov Random Fields: Theory and Applications*. Boca Raton: Chapman & Hall.
- Rue, H., Martino, S., and Chopin, N. 2009. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B* 71: 319–392.
- Sang, H. and Huang, J. 2012. A full scale approximation of covariance functions for large spatial data sets. *Journal of the Royal Statistical Society: Series B* 74: 111–132.
- Silverman, B. 1984. Spline smoothing: the equivalent variable kernel method. *Annals of Statistics*, 12: 898–916.
- Sollich, P. and Williams, C. 2005. Using the equivalent kernel to understand Gaussian Process regression. In *Advances in Neural Information Processing Systems 17*, MIT Press, pp. 1313–1320.
- Song, H., Fuentes, M., and Ghosh, S. 2008. A comparative study of Gaussian geostatistical models and Gaussian Markov random field models. *Journal of Multivariate Analysis* 99: 1681–1697.
- Stein, M. 1999. *Interpolation of Spatial Data : Some Theory for Kriging*. N.Y.: Springer.
- Stein, M. and Fang, D. 1997. Discussion of ozone exposure and population density in Harris County, Texas, by R.J. Carroll et al. *Journal of the American Statistical Association* 92: 408–411.

- Stein, M. L., Chi, Z., and Welty, L. J. 2004. Approximating likelihoods for large spatial data sets. *Journal of the Royal Statistical Society, Series B* 66: 275–296.
- Wall, M. 2004. A close look at the spatial structure implied by the CAR and SAR models. *Journal of Statistical Planning and Inference* 121: 311–324.
- White, G. and Ghosh, S. 2009. A stochastic neighborhood conditional autoregressive model for spatial data. *Computational Statistics & Data Analysis* 53: 3033–3046.
- Wikle, C. 2002. Spatial modeling of count data: A case study in modelling breeding bird survey data on large spatial domains. In *Spatial Cluster Modelling*, eds. A. Lawson and D. Denison, Chapman & Hall, pp. 199–209.
- Wolfinger, R. and O’Connell, M. 1993. Generalized linear mixed models: A pseudo-likelihood approach. *Journal of Statistical Computation and Simulation* 48: 233–243.
- Yue, Y. and Speckman, P. 2010. Nonstationary spatial Gaussian Markov random fields. *Journal of Computational and Graphical Statistics* 19: 96–116.
- Zhu, J., Huang, H., and Reyes, P. 2010. On selection of spatial linear models for lattice data. *Journal of the Royal Statistical Society: Series B* 72: 389–402.